# Computing at CDF

## Frank Wurthwein
*MIT/FNAL-CD*
for the CDF Collaboration

➢ **Introduction**

➢ **Computing requirements**
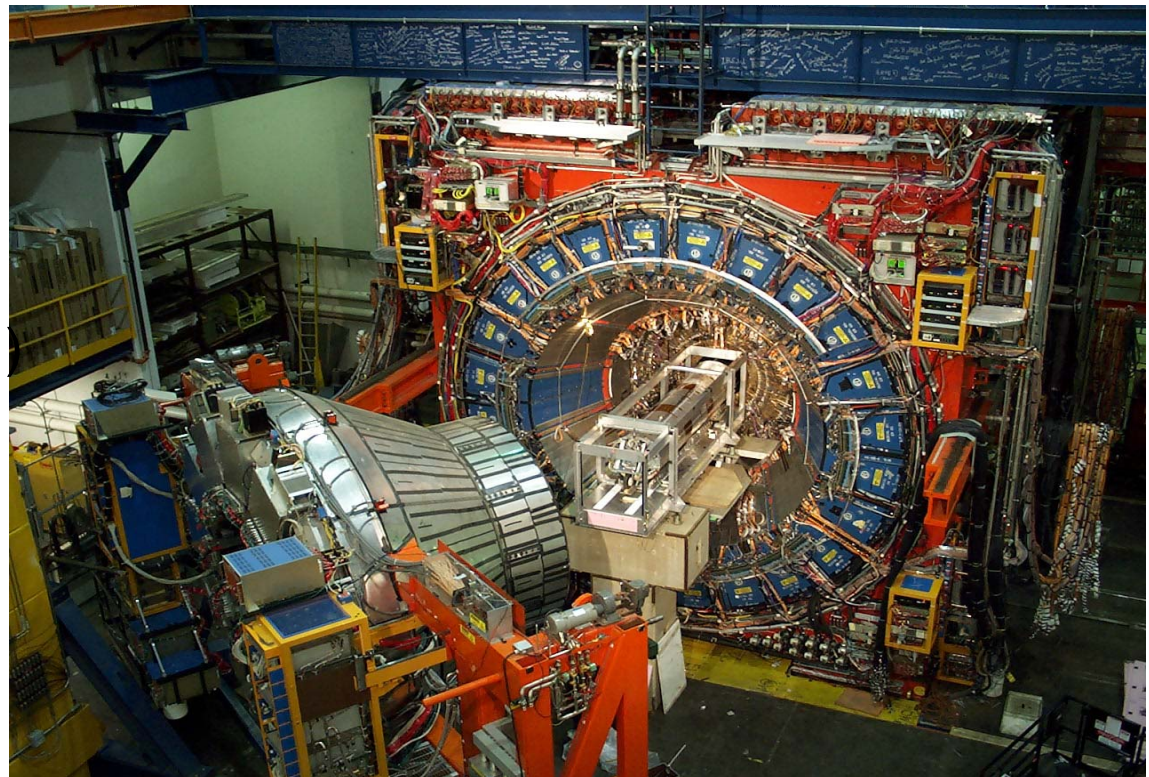
➢ **Central Analysis Farm**

➢ **Conclusions**

# CDF in a Nutshell

- CDF + D0 experiments analyze $p\bar{p}$ collisions from Tevatron at Fermilab
- Tevatron highest energy collider in world ($\sqrt{s} = 2$ TeV) until LHC
- Run I (1992-1996) huge success $\rightarrow$ 200+ papers (t quark discovery, ...)
- Run II (March 2001-) upgrades for luminosity ($\times 10$) + energy (~10%$\uparrow$)
  $\rightarrow$ expect integrated luminosity $20\times$ (Run IIa) and $150\times$ (Run IIb) of Run I

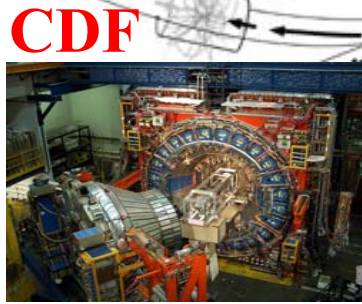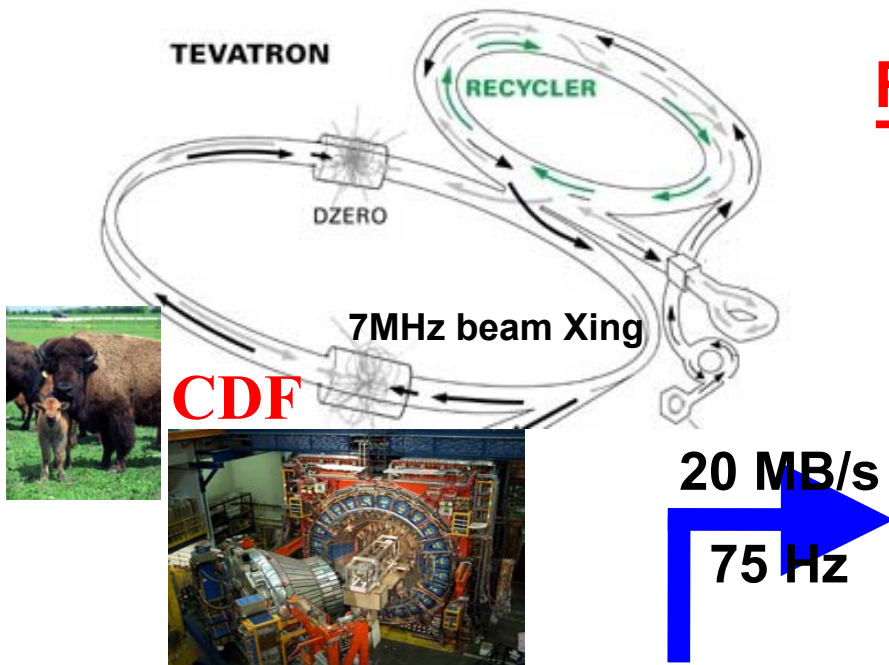## Run II physics goals:

- Search for Higgs boson
- Top quark properties ($m_t$, $\sigma_{tot}$, ...)
- Electroweak ($m_W$, $\Gamma_W$, $ZZ\gamma$, ...)
- Search for new physics (e.g. SUSY)
- QCD at large $Q^2$ (jets, $\alpha_s$, ...)
- CKM tests in *b* hadron decays

# CDF RunII Collaboration



525+ scientists
55+ institutions
11+ countries

Students
Postdoc's
Professors
Research Scientists

**Goal**: **Provide computing resources for 200+ collaborators simultaneously doing analysis per day!**

*Frank Wurthwein/MIT*

# CDF DAQ/Analysis Flow

User Desktops

TEVATRON

RECYCLER

DZERO

Robotic Tape Storage

7MHz beam Xing

Data    Analysis

CDF

20 MB/s

75 Hz

Read/write

Data

0.75 Million channels

L1
↓
L2
↓
300 Hz

Level 3 Trigger

MC    Recon

Production Farm

Central Analysis Facility (CAF)
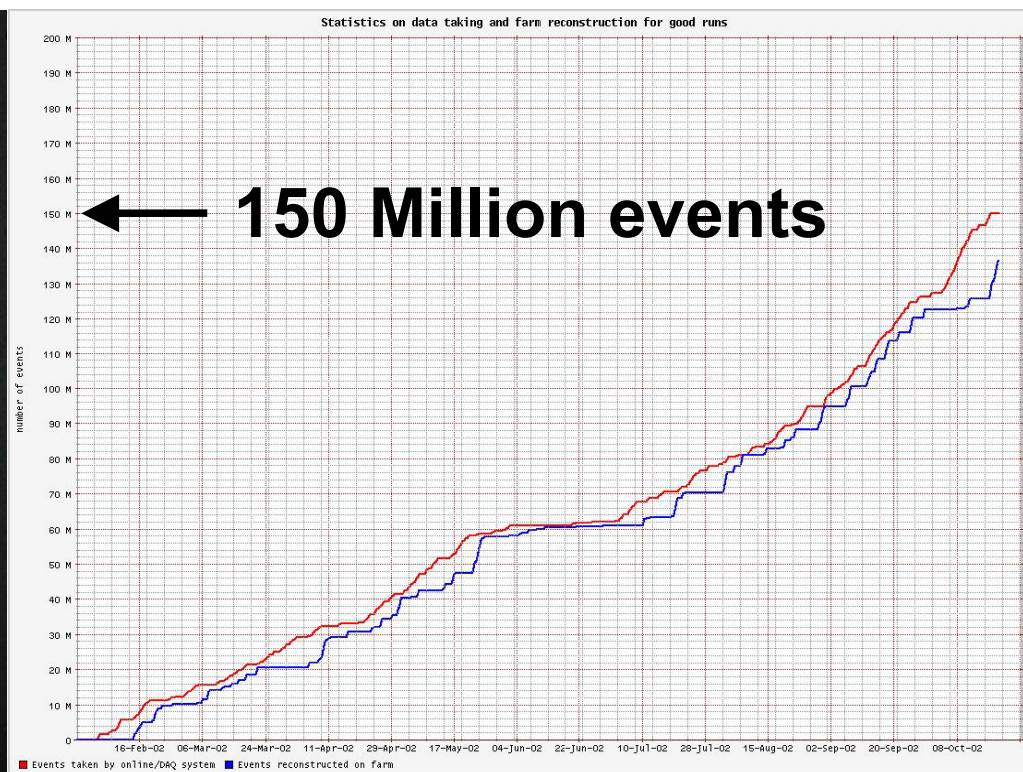
*Frank Wurthwein/MIT*

*LCCWS'02*

# Reconstruction Farms

**Data reconstruction + validation, Monte Carlo generation**

**154 dual P3's (equivalent to 244 1 Ghz machines)**

**Job management:**

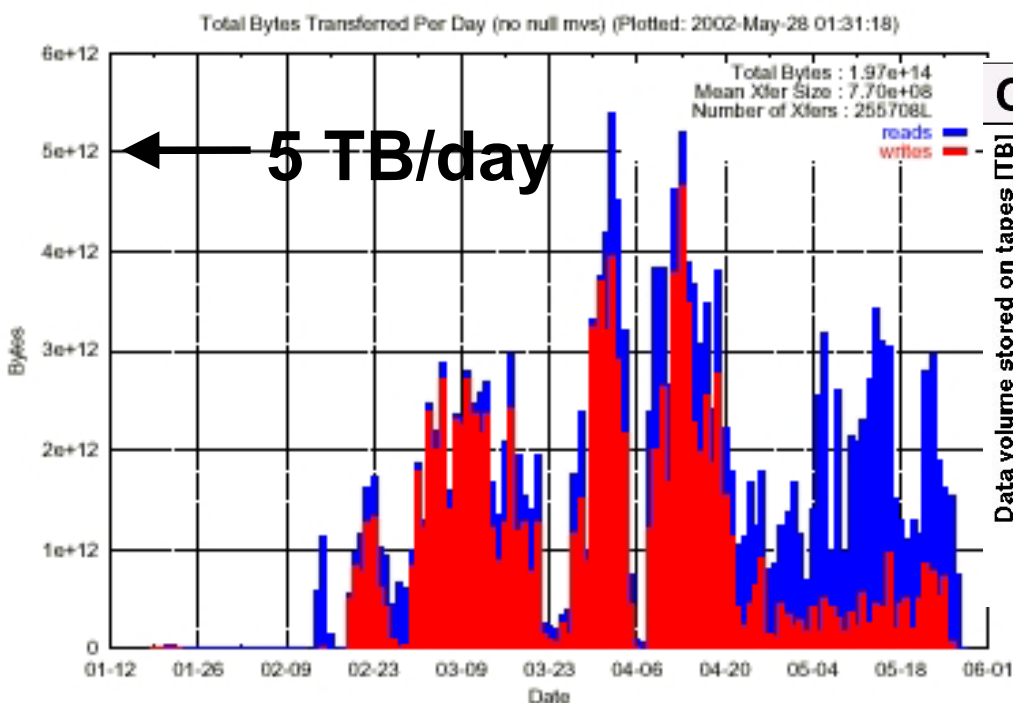➢ **Batch system → FBSNG developed at FNAL**
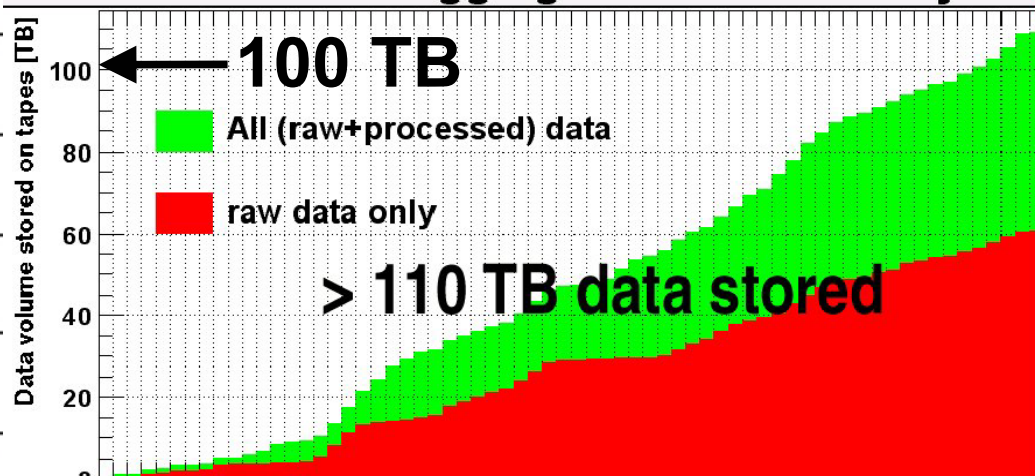
➢ **Single executable, validated offline**



Statistics on data taking and farm reconstruction for good runs

← **150 Million events**

Events taken by online/DAQ system    Events reconstructed on farm

# Data Handling

**Data archived using STK 9940 drives and tape robot**

**Enstore:** **Network-attached tape system developed at FNAL**
**→ provides interface layer for staging data from tape**

# Database Usage at CDF

## Oracle DB: Metadata + Calibrations

### DB Hardware:

- 2 Sun E4500 Duals

- 1 Linux Quad

### Presently evaluating:

- MySQL

- Replication to remote sites

- Oracle9 streams, failover, load balance



fcdfora1

fcdfora2

# **Data/Software Characteristics**

## Data Characteristics:
- **Root I/O sequential for raw data: ~250 kB/event**
- **Root I/O multi-branch for reco data: 50-100 kB/event**
- **'Standard' ntuple: 5-10 kB/event**
- **Typical RunIIa secondary dataset size: $10^7$ events**

## Analysis Software:
- **Typical analysis jobs run @ 5 Hz on 1 GHz P3**
   - $\rightarrow$ **few MB/sec**
- **CPU rather than I/O bound (FastEthernet)**

# Computing Requirements



| Fiscal Year | Lum (fb⁻¹) 4.1 | Batch CPU (THz) 4.7 | Farm CPU (THz) 1.3 | Static Disk (TB) 540 | Read Cache (TB) | Write Cache (TB) | Disk I/O (GB/s) 4.9 | Archive I/O (GB/s) | Archive Volume (PB) 1.7 |

**Requirements set by goal:**

200 simultaneous users to analyze secondary data set ($10^7$ evts) in a da...
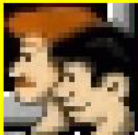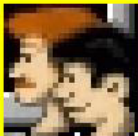
**Need ~700 TB of disk and ~5 THz of CPU by end of FY'05:**

→ need lots of disk→ need cheap disk → IDE Raid

→ need lots of CPU→ commodity CPU → dual Intel/AMD

# Past CAF Computing Model

Large SMP (128 processor SGI)
Expensive disks (FiberChannel/SCSI)

Analysis Code Development
Analysis Job Debugging
Interactive Analysis Jobs
Batch Jobs
"Other" Usage

fcdfsgi2

**Very expensive to expand and maintain**

**Bottom line:**
**Not enough 'bang for the buck'**

# Design Challenges

develop/debug interactively @ remote desktop
✓ code management & rootd

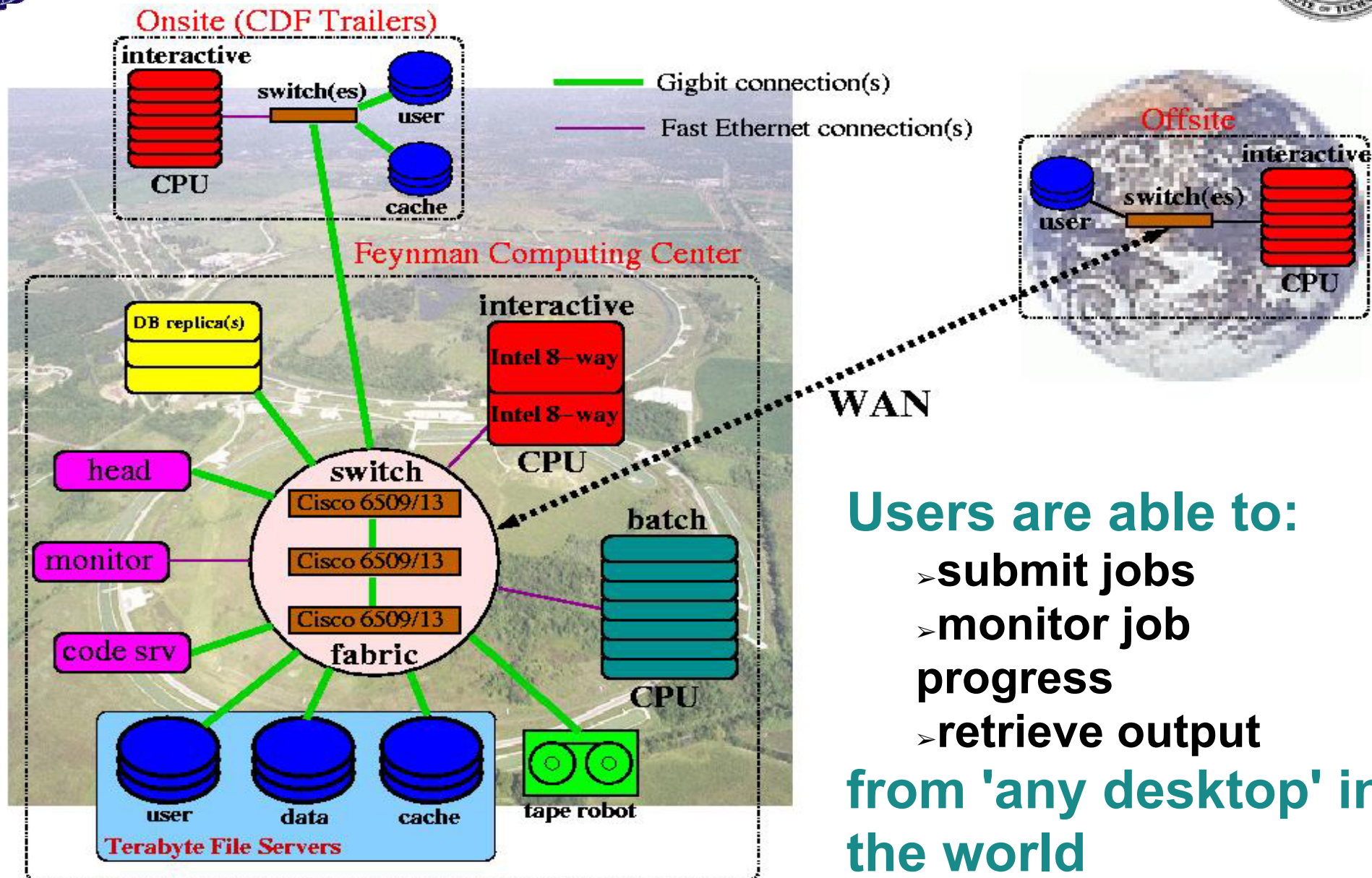Send binary & 'sandbox' for execution on CAF
✓ kerberized gatekeeper

no user accounts on cluster
BUT
user access to scratch space with quotas
✓ creative use of kerberos

# CAF Architecture



Users are able to:
- submit jobs
- monitor job progress
- retrieve output

from 'any desktop' in the world

# CAF Milestones

➤ **Start of CAF design** — 11/01

➤ **CAF prototype (protoCAF) assembled** — 2/25/02

➤ **Fully-functional prototype system (>99% job success)** — 3/6/02

➤ **ProtoCAF integrated into Stage1 system** — 4/25/02

➤ **Production Stage1 CAF for collaboration** — 5/30/02

**Design → Production system in 6 months!**



ProtoCAF

Stage1

# CAF Stage 1 Hardware



Code Server

File Servers

Worker Nodes

Linux 8-ways
(interactive)

# Stage 1 Hardware: Workers



**Workers** (**132 CPUs**, 1U+2U rackmount):

16 2U Dual Athelon 1.6GHz / 512MB RAM

50 1U/2U Dual P3 1.26GHz / 2GB RAM FE (11 MB/s) / 80GB job scratch each

# Stage 1 Hardware: Servers

**Servers (35TB total, 16 4U rackmount):**

2.2TB useable IDE RAID50 hot-swap
Dual P3 1.4GHz / 2GB RAM
SysKonnect 9843 Gigabt Ethernet card

# File Server Performance



**Local disk reads**

- 32k dd blocks
- 64k dd blocks
- 128k dd blocks
- 512k dd blocks
- 1M dd blocks
- 2M dd blocks

200 MB/s

60 MB/s

**Remote reads from CAF file server**

- NFS dd
- NFS Edm_ReadWriteSeqRoot
- Rootd Edm_ReadWriteSeqRoot
- NFS AC++ job
- Rootd AC++ job

70 MB/s

**Server/Client Performance**: Up to **200MB/s local reads, 70 MB/s NFS**

**Data Integrity tests**: md5sum of local reads/writes under heavy load
BER read/write = $1.1 +- 0.8 \times 10^{-15}$ / $1.0 +- 0.3 \times 10^{-13}$

**Cooling tests**: Temp profile of disks w/ IR gun after extended disk thrashing

# Stage2 Hardware

## Worker nodes:

238 Dual Athlon MP2000+, 1U rackmount

### 1 THz of CPU power

## File servers:

76 systems, 4U rackmount, dual red. Power supply
14 WD180GB in 2 RAID5 on 3ware 7500-8
2 WD40GB in RAID1 on 3ware 7000-2
1 GigE Syskonnect 9843
Dual P3 1.4GHz

### 150 TB disk cache

# Stage1 Data Access

**Static files on disk:**

> NFS mounted to worker nodes
> remote file access via rootd

**Dynamic disk cache:**

> dCache in front of Enstore robot

# Problems & Issues

## Resource overloading:

> DB meltdown $\rightarrow$ dedicated replica, startup delays
>
> Rcp overload $\rightarrow$ replaced with fcp
>
> Rootd overload $\rightarrow$ replaced with NFS,dCache
>
> File server overload $\rightarrow$ scatter data randomly

## System issues:

> Memory problems $\rightarrow$ improved burn-in for next time
>
> Bit error during rcp $\rightarrow$ checksum after copy
>
> dCache filesystem issues $\rightarrow$ xfs & direct I/O

# Lessons Learned

➢ **Expertise in FNAL-CD is essential.**

➢ **Well organized code management is crucial.**

➢ **Independent commissioning of data handling and job processing $\rightarrow$ 3 ways of getting data to application.**

# CAF: User Perspective

## Job Related:

- Submit jobs
- Check progress of job
- Kill a job

## Remote file system access:

- 'ls' in job's 'relative path'
- 'ls' in a CAF node's absolute path
- tail' of any file in job's 'relative path'

# CAF Software

# CAF User Interface

➢ **Compile, build, debug analysis job on 'desktop'**

➢ **Fill in appropriate fields & submit job**

section integer range

output destination

user exe+tcl directory

➢ **Retrieve output using kerberized FTP tools ... or write output directly to 'desktop'!**

# Web Monitoring of User Queues

**Each user a different queue**

**Process type for job length**

    **test**:      5 mins
    **short**:       2 hrs

    **medium**:  6 hrs
    **long**:       2 days

**This example:**
    **1 job → 11 sections**

**(+ 1 additional section automatic for job cleanup)**



Netscape: FBSWWW CAF list of queues

File   Edit   View   Go   Communicator     Help

**FBSNG on the web**
Farm:         CAF
Time:        Thu May 23 02:32:41 2002
Report:     List of queues

Queues  Jobs  Nodes  Process Types

User Monitor

| Name | Status | Default Process Type | Share | Prio | Waiting | Ready | Running | Total |
|------|--------|----------------------|-------|------|---------|-------|---------|-------|
| akorn | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| amitl | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| anikeev | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| belforte | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msmartin | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| msn | OK | short | 1.00 | 0 | 1 | 0 | 11 | 12 |
| pauly | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| paus | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| ratnikov | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| rescigno | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| semeria | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sfiligoi | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sgromoll | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| shepard | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| sidoti | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| spezziga | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| test | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thkim | OK | short | 1.00 | 0 | 0 | 0 | 0 | 0 |
| thom | OK | short | 1.00 | 0 | 1 | 0 | 1 | 2 |

*LCCWS'02*

# Monitoring jobs in your queue

**Monitoring sections of your job**

# CAF Utilization



| ptype | Average % |
|-------|-----------|
| Short | 59.1 |
| Medium | 18.9 |
| Long | 12.9 |
| All processes | 91.0 |

Created on Oct 20 22:02:05 2002

## Active queues (last 24h)



Updated: Jun 18 12:20:03 2002

schuster · test · fkw · anikeev · tomohiro · gotra · akorn · castro
casarsa · cjl · jdlee · zyhanv · ikfuric · lys · nigmanov · gpope
jmuelmen · rescigno · dbstarr · ikrav · daronco · thom

Built using RRDTool

## Summary Table

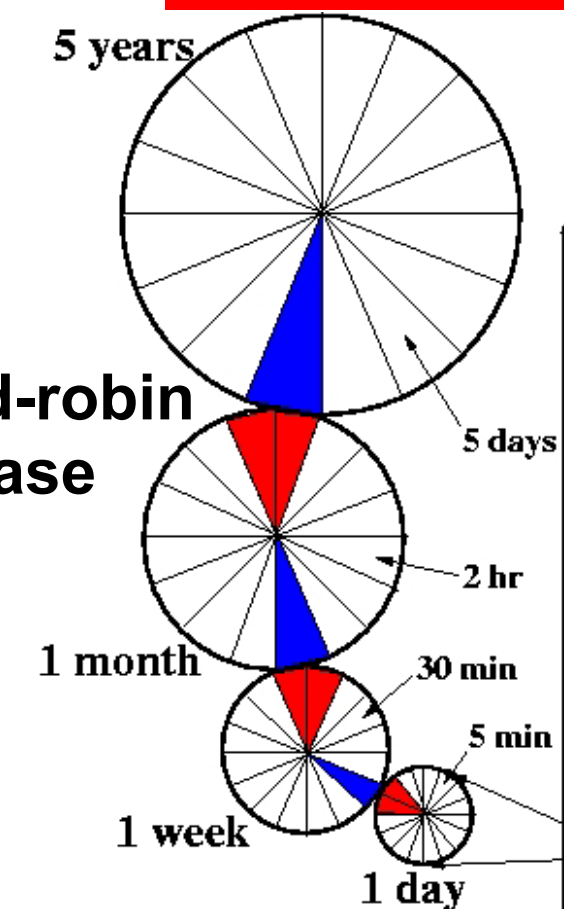| | Short | Medium | Long | All Types |
|---|---|---|---|---|
| Running sections | 98 | 7 | 21 | 126 |
| Pending sections | 0 | 70 | 66 | 136 |
| Waiting time [hh:mm] (24h average): | | | | |
|   per job | 0:52 | 0:15 | 0:00 | 0:33 |
|   per section | 4:14 | 3:29 | 2:44 | 3:29 |
| Running time [hh:mm] (24h average) | 0:27 | 4:34 | 4:08 | 3:03 |

Updated: Oct 20 22:00:03 2002

**CAF in active use by CDF collaboration**

➢ **300 CAF Users (queues) to date**
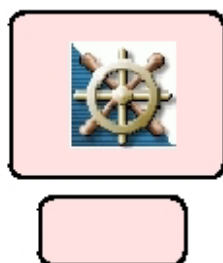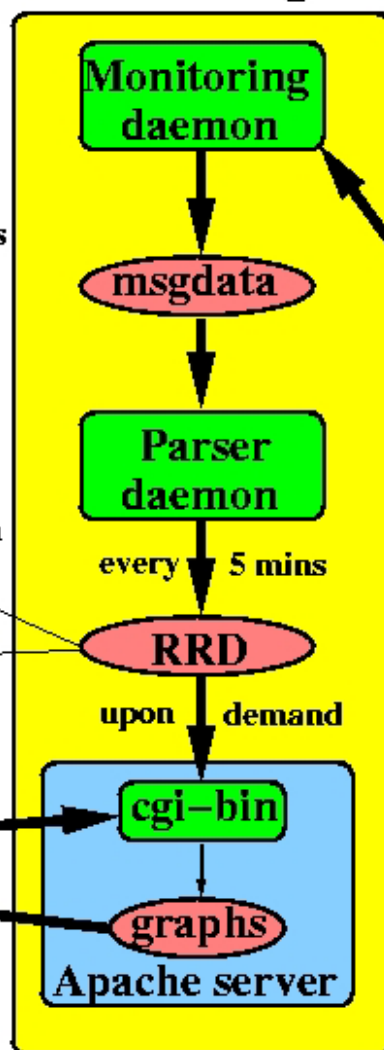➢ **Several dozen simultaneous users in a typical 24 hr period**

# CAF System Monitoring

# CPU Utilization



**CAF utilization steadily rising since opened to collaboration**

**Provided 10-fold increase in analysis resources for last summer physics conferences**

**Need for more CPU for winter**
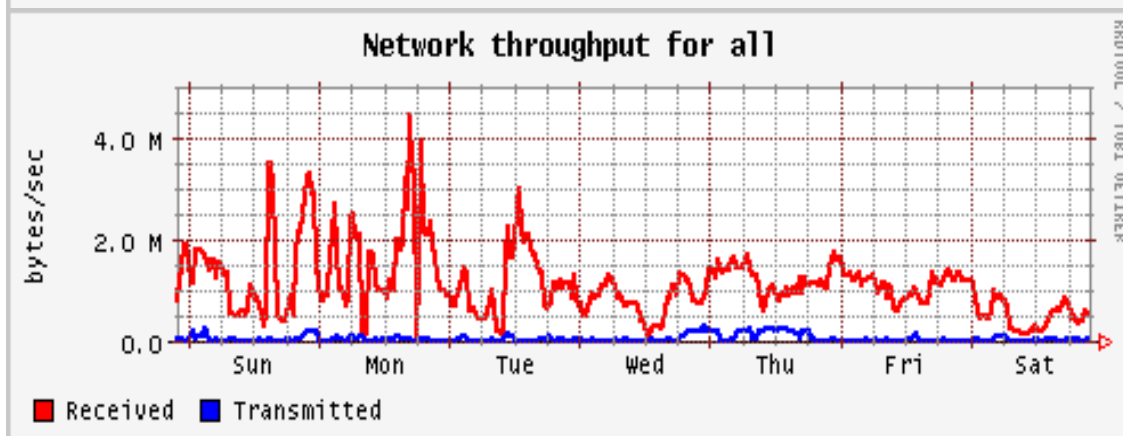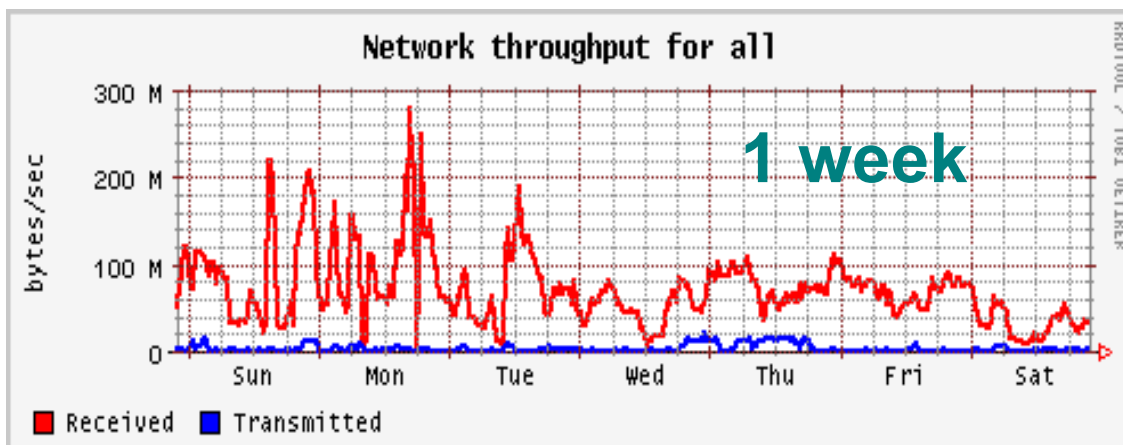
*LCCWS'02*

# Data Processing

## File Server

**Aggregate I/O**
**4-8TB/day**

**Aggregate I/O**

## Worker Node

**Average I/O**
**1-2MB/sec @**
**~80%CPU util.**

*LCCWS'02*

# Work in Progress

Stage2 upgrade: 1THz CPU & 150TB disk

SAM $\rightarrow$ framework for global data handling/distribution

"DCAF" $\rightarrow$ remote "replicas" of CAF

Central login pool @ FNAL

# CAF Summary

**Distributed Desk-to-Farm Computing Model**

**Production system under heavy use:**

- **Single farm at FNAL**
  - 4-8TB/day processed by user applications
  - Average CPU utilization of 80%
- **Many users all over the world**
  - 300 total users
  - typical: 30 users per day share 130 CPU's
  - Regularly several 1000 jobs queued
- **Connected to tape via large cache**
- **Currently updating to 1THz & 150TB**

# CDF Summary

**Variety of computing systems deployed:**

- ➢ **Single app. Farms: Online & Offline**

- ➢ **Multiple app. Farm: user analysis farm**

- ➢ **Expecting 1.7Petabyte tape archive by FY05**

- ➢ **Expecting 700TB disk cache by FY05**

- ➢ **Expecting 5THz of CPU by FY05**

- ➢ **Oracle DB cluster with loadavg & failover for metadata.**

*LCCWS'02*